

ROFusion: Efficient Object Detection using Hybrid Point-wise Radar-Optical Fusion

Liu Liu, Shuaifeng Zhi*, Zhenhua Du, Li Liu, Xinyu Zhang, Kai Huo, and Weidong Jiang

College of Electronic Science, National University of Defense Technology
liuliucn@outlook.com, zhishuaifeng@outlook.com

Abstract. Radars, due to their robustness to adverse weather conditions and ability to measure object motions, have served in autonomous driving and intelligent agents for years. However, Radar-based perception suffers from its unintuitive sensing data, which lack of semantic and structural information of scenes. To tackle this problem, camera and Radar sensor fusion has been investigated as a trending strategy with low cost, high reliability and strong maintenance. While most recent works explore how to explore Radar point clouds and images, rich contextual information within Radar observation are discarded. In this paper, we propose a hybrid point-wise Radar-Optical fusion approach for object detection in autonomous driving scenarios. The framework benefits from dense contextual information from both the range-doppler spectrum and images which are integrated to learn a multi-modal feature representation. Furthermore, we propose a novel local coordinate formulation, tackling the object detection task in an object-centric coordinate. Extensive results show that with the information gained from optical images, we could achieve leading performance in object detection (97.69% recall) compared to recent state-of-the-art methods FFT-RadNet [17] (82.86% recall). Ablation studies verify the key design choices and practicability of our approach given machine generated imperfect detections. The code will be available at <https://github.com/LiuLiu-55/ROFusion>.

Keywords: Radar-Optical Fusion · Object Detection · Deep Learning.

1 Introduction

Autonomous driving and Advanced Driver Assistance Systems (ADAS) often rely on different types of sensors to acquire a reliable perception. Mainstream sensors equipped in automotive vehicles are camera, Lidar and Radar, which are fused together thanks to their unique working mechanism and specialties. Existing mainstream multi-sensor fusion strategy uses camera and Lidar sensors for 3D object detection [2,21]. Mainly because Lidar owns a high angular resolution and range detection accuracy in a way of dense point clouds, and is complementary to camera images which are rich in contextual and semantic information of

* Shuaifeng Zhi is the corresponding author.

scenes. However, both camera and Lidar suffer from huge performance degradation in adverse weather conditions, which is a crucial requirement for long-term stable autonomous driving.

Radars are active sensors that measure the environment from reflected electromagnetic waves. Compared to Lidar, Radar has a robust capacity in severe weather conditions and can detect objects and obstacles within 250m with their distances and velocities. Furthermore, its low deployment cost makes Radar a requisite sensor in assistance systems. Radar data have developed different types of representations, including Radar occupancy grid maps, micro-Doppler signature, dense Range-Doppler-Azimuth (RAD) tensors and point clouds, with various processing costs and representational capacity.

Despite Radar’s advantages in stable and long-term scene perception, there have been few investigations on fusing Radar with other sensors in this task. This is partly caused by its entirely different imaging mechanism in contrast to cameras and Lidars, leading to extremely sparse point clouds or intuitive dense RAD spectrum, and relatively low elevation angular resolution as well. Fortunately, this problem has been partly solved with the development of the 4D imaging Radar, with a high angular resolution of about 1° in both azimuth and elevation. Some recent works also tried to conduct image-Radar fusion to alleviate the high sparsity of Radar point clouds [14,10,9].

Motivated by the above-mentioned challenges, we propose ROFusion, a hybrid point-wise approach to fuse Radar and camera data. Different from previous work in Radar-optical fusion, we seek to fuse dense contextual features from both modalities. We first acquire Radar and camera features respectively from single-modality extractors [17,8], and then use image-Radar association and hybrid point-based fusion strategy to merge cross-modality features at multiple hierarchies. Finally, a local coordinate formulation is proposed to decompose our tasks into classification and regression in an object-centric manner. Our method achieves a new state-of-the-art performance in both easy and hard cases of public RADial dataset [17].

To summarize, our contributions are as follows:

- We propose a hybrid point-wise fusion strategy to effectively associate dense Radar and image features.
- We propose a local coordinate formulation that simplifies object detection by classification and regression sub-tasks in an object-centric manner.
- We conduct extensive experiments on the RADial [17] benchmark and achieve a new state-of-the-art detection performance, with a significant boost over Radar-based baseline.

2 Related work

2.1 Point-based Methods

PointNet [15] designs a novel type of neural network that directly consumes the point cloud, which makes point-based detection methods process. For Radar

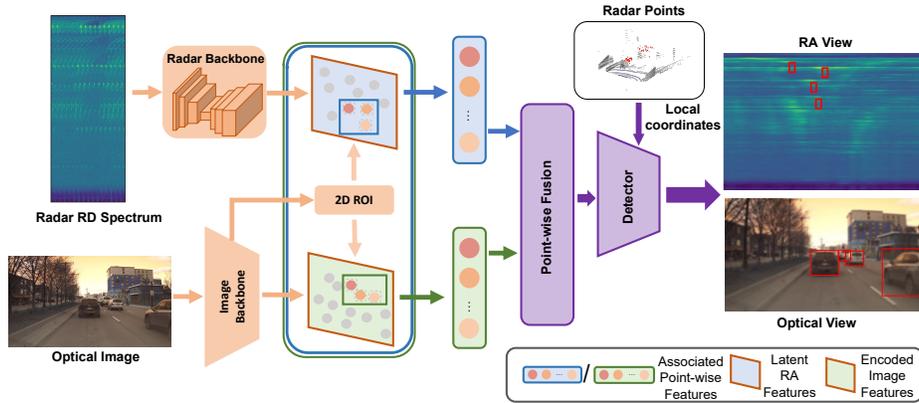


Fig. 1. ROFusion network architecture. Latent RA feature maps and camera-encoded semantics are first obtained by Radar and image backbone. The image 2D bounding boxes are used to associate the image and Radar feature maps via point cloud candidates. The point-wise module is next used for feature-level fusion, proposing a hybrid point feature to produce final Radar detection formulated in object-centric local coordinates.

point clouds, sparse structures take a challenge to object detection. One strategy [3] is to accumulate radar points into a dense occupancy grid mapping (OGM). For lightweight demand, [19] utilize novel point structure [16]. With the sparsity issue, [13] observes that a global message could enhance perception performance.

2.2 Camera-Radar Fusion Methods

Complementary information gives the opportunity for sensor fusion between the camera and Radar. Radar extracts the distance and velocity of objects, while semantic information is captured by cameras. There are normally three fusion levels between Radar and camera: early level, feature level, and late level. Radars are often used to generate the region of interest (ROI) for early-level fusion. Then, the predicted region is processed as an auxiliary refining optical task [4,7], which is computationally expensive. The decision level contrary utilizes two sources independently detect, proposing a strategy [1,24] defining whether one of the sensors failed. With different probability spaces, late-level fusion could not efficiently exert the capability of two sensors.

A naive approach is fusing Radar and camera in a latent feature space where the key point is Radar-camera association. CramNet [9] applies a dynamic voxelization fusing Radar and camera features, projecting each camera pixel with a 3D ray to rectify its location, which makes a robust performance for sensor failure. In [14], authors propose a frustum association that fully exploits Radar vertical information. CRAFT [10] also associates Radar and image, but implements them in a polar coordinate to handle the discrepancy between the coordinate system and spatial properties. The feature maps are then fused by a consecutive cross-attention strategy.

3 Method

In this section, we present architectural details of our method ROFusion as well as key design choices enabling accurate object detection with a hybrid point-wise optical-Radar fusion. An overall architecture is provided in Figure 1. We first take Radar RD spectrums with corresponding images as our network inputs and extract their dense features. Radar points filtered by prior information such as image detection bounding boxes are then adopted as anchors to associate Radar and image features. Furthermore, a hybrid point-wise fusion complements the surrounding semantics of targets to produce new point features. A detection header finally predicts object locations in per-object local coordinates.

To summarize, our pipeline consists of three main modules including: (1) a dense feature extraction module from both RGB images and corresponding high-definition RD tensors to acquire contextual information of scenes (Section 3.1); (2) a hybrid point-based fusion module to associate dense Radar embedding of scattering points with image features (Section 3.2); (3) a local coordinate module formulating object detection task in an object-centric manner (Section 3.3). We finally show the training configurations of our method in Section 3.4.

3.1 Dense Feature Extraction

In order to acquire rich contextual information about objects within 3D scenes, we leverage dense convolutional neural networks (CNNs) to extract dense feature embedding of both Range-Doppler (RD) Radar maps and camera image observations.

Radar Feature Extractor Radar-based scene understanding from its Range-Doppler (RD) map has recently gained attention as it contains all information on range, azimuth and elevation. In addition, the RD map owns less computational acquisition costs and is a dense representation compared to Range-Doppler-Azimuth (RAD) tensors and sparse point clouds, respectively. We propose to use a dense CNN model as our Radar backbone module, inspired by FFT-RadNet [17]. Specifically, it aims to learn a multi-scale dense representation of Range-Azimuth (RA) maps from their input RD counterparts, with a tailored MIMO pre-encoder [17]. In this manner, we seek to learn a dense feature embedding of RA maps as they are closely related to downstream vehicle detection tasks.

Image Feature Extractor To enrich radar features with optical image features, we encode the corresponding RGB image into a dense feature embedding with a vision CNN model. To reduce computational overhead, we simply use an ImageNet [18] pre-trained ResNet-18 model [8] and keep the weights intact during training. Please note that our image backbone module could be replaced by stronger vision models such as ResNet-152 [8] and vision Transformers [20], depending on the computation budgets.

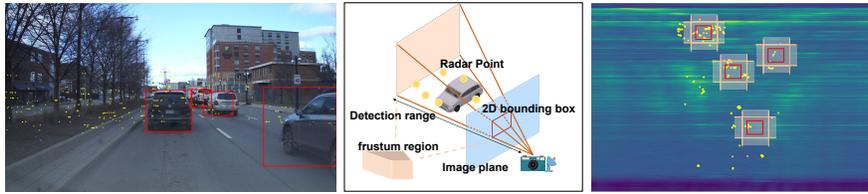


Fig. 2. Image-Radar Association. 2D Detector using image features provides the azimuth of interest (left), which leads to a frustum region to select candidate object-related point clouds (middle). We filter noise and background points depending on whether their relative radial or angular distance to the object center is beyond a certain threshold or not, as discussed in Section 3.3 (right).

3.2 Point Fusion

Image-Radar-Association In this section, we explain how to establish a cross-modality association of target objects with provided sensor calibration information and prior optical detection results.

As dense features of RA maps and images are difficult to conduct dense alignment due to their different imaging mechanism, we rely on Radar point clouds to bridge them at point-level. Specifically, we represent each Radar point as a 3D point $p = (r, a, d, u, v, x, y, z)$, where (r, a, d) and (x, y, z) are its locations within RAD tensors and real world coordinates, respectively. With the intrinsic and extrinsic of the camera model, we transform Radar points into image coordinates as follows:

$$u = f_x \frac{x'}{z'} - p_x, \quad v = f_y \frac{y'}{z'} - p_y, \quad (1)$$

where (f_x, f_y, p_x, p_y) are camera intrinsic parameters, (x', y', z') is 3D position within camera coordinate transformed by camera extrinsic $[R|t]$ and (x, y, z) .

2D object bounding boxes within images are treated as Region of Interest (ROI) filters separating the region of interests out of background and noises, as explained in Figure 2. The 2D bounding boxes provide strong prior angular information of objects, eliminating the uncertainty caused by Radar sidelobe jamming. At this stage, we treat all points within these 2D ROIs as candidate points for the next fusion stage. However, these boxes within images cannot cope with range estimation, as points within 3D space in the cone area (middle of Figure 2) all project within the 2D ROIs. To address the range inaccuracy, we consider a local coordinate strategy as detailed in Section 3.3.

Hybrid Point Fusion We propose a point-based method that generates per Radar point fused feature from pixel-level RA and image features. Inspired by DenseFusion [22], we implement a variant architecture that fuses semantics and velocity.

Assume there are k 3D Radar points from the previous association stage, we collect pixel-wise features from encoded RA features F_R and semantic image features F_I , respectively. Concretely, with a set of k point clouds $P = \{p_1, p_1, \dots, p_k\}$, we extract corresponding per-pixel features $F_R = \{F_r^{p_1}, F_r^{p_2}, \dots, F_r^{p_k}\}$ and RA features $F_I = \{F_i^{p_1}, F_i^{p_2}, \dots, F_i^{p_k}\}$, where $F_r^{p_k}$ and $F_i^{p_k}$ are pixel-level features

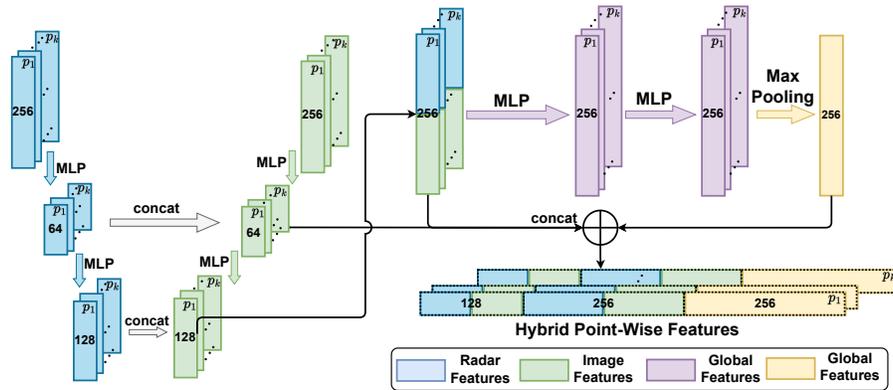


Fig. 3. Hybrid point fusion architecture. Extracted dense Radar and image features are processed by a series of MLPs and fused at multiple scales. With a per-object global feature from max-pooling across points, we reach a hybrid multi-modal feature with spatial and semantic information by concatenating all fused features of various scales.

from RA and image of point p_k . As shown in Figure 3, considering the difference within local distribution and semantics of these two feature spaces, the obtained point-wise features are combined in a hierarchical manner. As low-level and high-level fusions are both efficiently discriminative point-level features, we fuse them at different scales via concatenation after being sequentially processed by a set of shared MLPs. Another key point here is to obtain a per-object global contextual feature which, in principle, reveals the attributes of the same target which shares across domains. The global point-level feature is obtained via a max-pooling operation of fused features across all candidate points of the same object. We obtain a set of hybrid point-wise features by concatenating all above mentioned fused features at various scales. These features are fed into a detector that predicts per-point object center locations (see Section 3.4).

3.3 Object-centric Local Coordinates

We have experimentally found that directly regressing object locations is not only challenging to achieve purely from extremely sparse point clouds, but also involved with the absolute scale of sensing environments, imaging resolutions and object locations. To tackle this problem, we propose to decompose object detection task into a combination of classification and regression sub-tasks at an object-centric local coordinate.

As shown in Figure 4, we establish a new coordinate whose origin is at object center, and x , y axes are parallel to range and azimuth axes of RA. For Radar points within 2D bounding boxes, we encode their relative distance to the center position of targets at both axes based on a set of discrete bins at a certain resolution. We further predict a residual offset via regression on top of classification results to reach the final localization prediction. The motivation comes from the fact that we only focus on the features around the target and this formulation

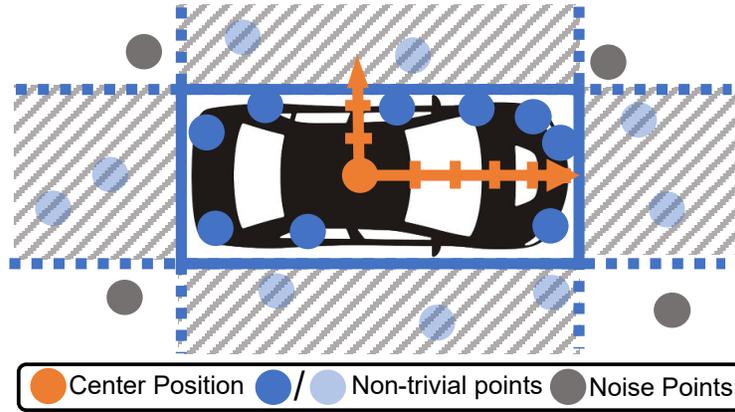


Fig. 4. Illustration of the local coordinates originated at object center position. Non-trivial points lying within or around the object are kept during training, while noise points are discarded.

decouples the network prediction from the above-mentioned imaging conditions. In this paper, the relative distance between center labels and Radar points is modeled by discretizing object dimension into 5 bins and 11 bins for azimuth and range, respectively.

It is also worth noting that although Radar points around objects are considered candidates, points which also lie within bounding boxes but are reflected by context near objects also have valuable information. We term such points ‘non-trivial’ points if their relative distance at *any* axes satisfies our above-mentioned discretization. For example, points reflected by non-object regions within boxes may have a large variation w.r.t. range dimension, but share a large correlation to object at angle dimension. In such cases, we may still include these points as training data and only penalize our network prediction by the deviation at the angle-axis prediction. This investigation also filters radar foreground and background points, eliminating range uncertainty as explained in Section 3.2.

In the training loop, we use non-trivial points as a data augmentation, which could partly relieve the spatial sparsity of object Radar points. All other points are regarded as background points or noises and are not involved during network training and inference.

3.4 Object Detection and Training Configurations

As described in Section 3.3, the detection task is divided into two parts, a RA map coarse classification and a refined regression. The two-part predictions are trained with a combined loss composed of a Cross-Entropy loss and a Smooth-L1 loss [6]. Denote the network prediction of classification and regression as $\hat{y}_{cls}^{B \times N \times 16}$ and $\hat{y}_{reg}^{B \times N \times 2}$, the training loss is:

$$\mathcal{L} = \mathcal{L}_{CrossEntropy}(y_{cls}, \hat{y}_{cls}) + \alpha \mathcal{L}_{Smooth-L1}(y_{reg}, \hat{y}_{reg}), \quad (2)$$

where $\alpha = 10$ is a weight balancing parameter.

In the training phase, we use the object’s 2D ground truth bounding box to get a precise association. In the test phase, a pertained 2D detector is used to provide object bounding boxes for evaluation. Specifically, we choose to use a pre-trained off-the-shelf YoloX [5] to obtain 2D bounding boxes on testing images without any fine-tuning on the target datasets. We also show the performance of our method given oracle bounding boxes as a limited case to show the potential upper bound performance of our method.

4 Experiments

4.1 Dataset and Metrics

Dataset We evaluate our model on the RADIAL dataset [17] consisting of RD spectrums and Radar points of a high-definition Radar with corresponding camera observations. Its 91 sequences are divided into hard and easy cases depending on the intensity of Radar perturbation. We strictly follow the official splits into training, validation, and test division at a portion of 70%, 15%, 15%. Since our proposed point-based architecture requires there are Radar reflection peaks from objects, we remove training candidates where no Radar points are included within object bounding boxes.

Metric The evaluation metrics for object detection are Average Precision (AP) and Average Recall (AR), given a validated positive prediction whose Intersection-over-Union (IoU) to the ground truth is greater than 50% [17]. We also present the absolute Range and Angle error to analyze the prediction accuracy.

4.2 Baseline

Implement Details We implement our image backbone with a pre-trained ResNet-18 [8] model. The color image is of size 960×540 and we use the semantic features of the last layer as dense image features. The Radar backbone adopts the design of FFT-RadNet [17] while we further simplify the FPN [12] model to reduce computational complexity. Due to the high definition nature of used Radar sensor, $\frac{1}{4}$ of native resolution is taken and has been proven to be enough for near-by object discrimination [17]. We train our ROFusion model for 40 epochs with a batch size of 8 and 1×10^{-4} learning rate with Adam optimizer [11] on a single NVIDIA Tesla V100 GPU. During inference, the bottom Radar points inside objects’ 2D bounding boxes are considered as sensor-facing endpoints and are used to generate the heuristic object-centric local coordinates (hLC).

Results In Tabel 1 we report object detection results of our method compared to leading state-of-the-art methods FFT-RadNet [17] and baseline method Pixor [23]. Ground truth bounding boxes and object positions are used to demonstrate the effectiveness of ROFusion. We also evaluate the performance of our method

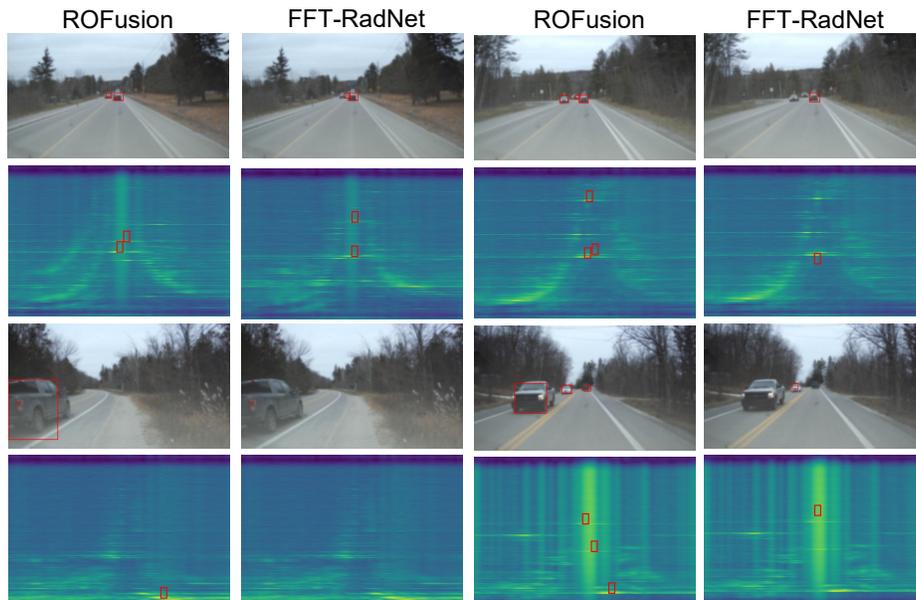


Fig. 5. Qualitative results for object detection from Camera (row 1 & 3) and RA (row 2 & 4) view. In the RA plots, the detection boxes are presented corresponding to RA dense maps in Euclidean space.

Table 1. Object detection performance on RADial dataset [17]. AR (%) is computed with an IoU threshold of 50%. $R(m)$ and $A(^{\circ})$ indicate the mean Range and Angle error. PC, RA, RD and IM mean point clouds, range-azimuth maps, range-doppler maps and images, respectively.

Method	Input	Overall			Easy			Hard		
		AR(%) \uparrow	$R(m)$ \downarrow	$A(^{\circ})$ \downarrow	AR(%) \uparrow	$R(m)$ \downarrow	$A(^{\circ})$ \downarrow	AR(%) \uparrow	$R(m)$ \uparrow	$A(^{\circ})$ \downarrow
Pixor [23]	PC	32.32	0.17	0.25	28.83	0.15	0.19	38.69	0.19	0.33
Pixor [23]	RA	81.68	0.10	0.20	88.02	0.09	0.16	70.10	0.12	0.27
FFT-RadNet [17]	RD	82.18	0.11	0.17	91.69	0.10	0.13	64.82	0.13	0.26
FFT-RadNet* [17]	RD	82.86	0.12	0.11	93.12	0.11	0.10	64.13	0.15	0.13
Ours	IM+RD+PC	97.69	0.12	0.21	97.79	0.11	0.19	97.52	0.12	0.22
Ours-hLC	IM+RD+PC	93.64	0.12	0.23	95.21	0.12	0.22	91.22	0.13	0.25

- We denote that FFT-RadNet* [17] as detector with 0.5 discrimination threshold for a fair comparison, using authors' provided weights.

with proposed heuristic local coordinates estimation for practical purposes. Although the radar sparsity causes the worse Angle error, we have observed a clear performance boost despite using sparse point-level features thanks to the optical information and our local coordinate formulation. Our method outperforms [17] at overall recall rate with a gap of +14.83%. It is worth to highlight that our hybrid point-wise fusion scheme achieves a promising +27.65% recall boost and a 0.12m Range error in the hard cases, overcoming interference problems caused by Radar noise to the dense formulation in [17]. Qualitative results can be found in Figure 5.

Table 2. 2D Object detection metrics of YOLOX Network [5] on the test set.

Method	Input	Overall		Easy		Hard	
		AP(%)	AR(%)	AP(%)	AR(%)	AP(%)	AR(%)
YOLOX [5]	IM	90.48	91.03	89.79	91.86	91.77	89.54

Table 3. Detection performance on RADIAL [17] given predicted 2D boxes.

Method	Overall				Easy				Hard			
	AP(%)	AR(%)	R(m)	A(°)	AP(%)	AR(%)	R(m)	A(°)	AP(%)	AR(%)	R(m)	A(°)
FFT-RadNet* [17]	97.39	82.86	0.12	0.11	98.96	93.12	0.11	0.10	93.46	64.13	0.15	0.13
Ours(YOLOX)-hLC	91.58	95.15	0.13	0.21	91.03	96.07	0.13	0.20	92.63	93.47	0.13	0.23

To further demonstrate the practicability of our method, we use network predicted detection results to conduct the evaluation. We first reveal the quality of the adopted YOLOX [5] 2D detector in Table 2, with a moderate performance drop in terms of optical detection accuracy, it is expected that the imperfect 2D detection results would affect the filtering process of our pipeline. Table 3 compares the *AP* and *AR* metrics with machine-generated bounding boxes in both easy and hard cases as well. While the *AP* metric lack behind baseline model due to the quality of network inferred 2D bounding boxes, our method with YOLOX [5] 2D detector still achieves a higher *AR* metric for both overall and especially difficult cases. The *AR* performance gain comes from the 2D bounding boxes association and heuristic local coordinates. The Range error of our method also outperforms FFT-RadNet* [17] in difficult cases. These results show that our local coordinates successfully extract the range information for sparse Radar points even though the prior 2D detection is less accurate.

4.3 Ablation

In this section, we conduct ablative experiments to validate the key components of our method: local coordinate (LC) formulation and image fusion module (IM).

As shown in Table 4, we have shown two variants of ROFusion: Ours (w/o LC) is the variant where we remove the local coordinate formulation in the training stage, but conduct the two classification and regression sub-tasks in the original RA maps; we also remove the image level feature fusion module out of training process. From the statistics, we conclude that the local coordinate formulation is significant in enabling accurate learning from spare Radar point clouds. The integration of image features gives further performance boost upon competing performance. It is also worth noting that in addition to the image feature, the prominent prior information introduced by optical detection is another key factor supporting the overall learning process from sparse Radar point clouds.

5 Conclusion

In this paper, we present ROFusion, a novel point-wise Radar-Optical fusion network for object detection. We have demonstrated that our method could effectively exploit camera semantics to enhance Radar detection. With hybrid point fusion and local coordinate formulation, ROFusion achieves state-of-the-art

Table 4. Ablation study on the key components of ROFusion.

Method	Input	Overall			Easy			Hard		
		AR(%) \uparrow	R(m) \downarrow	A($^\circ$) \downarrow	AR(%) \uparrow	R(m) \downarrow	A($^\circ$) \downarrow	AR(%) \uparrow	R(m) \downarrow	A($^\circ$) \downarrow
Ours (w/o LC)	RD+PC	27.4	0.22	0.46	32.94	0.24	0.40	15.78	0.20	0.57
Ours	RD+PC	96.58	0.09	0.22	96.31	0.08	0.22	96.85	0.09	0.24
Ours	IM+RD+PC	97.69	0.12	0.21	97.79	0.11	0.19	97.52	0.12	0.22

performance on the public RADial dataset [17], showing the potential capability for multi-sensor fusion. However, our method still relied on the quality of 2D object detection as prior information to filter potential object Radar points. In addition, considering the difference in imaging mechanism, more in-depth analysis of camera-Radar fusion stratify at the feature level is worth investigating, possibly aided by a powerful Transformer backbone using attention mechanism. This can be an exciting venue for our future work.

Acknowledgements

This work was partially supported by the National Key Research and Development Program of China (Grant No. 2021YFB3100800), the National Natural Science Foundation of China (Grant No. 61921001, 62201603) and Research Program of National University of Defense Technology (Grant No. ZK22-04).

References

1. Josip Ćesić, Ivan Marković, Igor Cvišić, and Ivan Petrović. Radar and stereo vision fusion for multitarget tracking on the special euclidean group. *Robotics and Autonomous Systems*, 83:338–348, 2016. **3**
2. Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1907–1915, 2017. **1**
3. Maria Dreher, Emeç Erçelik, Timo Bänziger, and Alois Knol. Radar-based 2d car detection using deep neural networks. In *Proceedings of the International IEEE Conference on Intelligent Transportation Systems (ITSC)*, pages 1–8, 2020. **3**
4. Floris Gaisser and Pieter P Jonker. Road user detection with convolutional neural networks: An application to the autonomous shuttle wepod. In *Journal of Machine Vision and Applications*, pages 101–104. IEEE, 2017. **3**
5. Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. YOLO: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430*, 2021. **8, 10**
6. Ross Girshick. Fast r-cnn. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 1440–1448, 2015. **7**
7. Xiao-peng Guo, Jin-song Du, Jie Gao, and Wei Wang. Pedestrian detection based on fusion of millimeter wave radar and vision. In *International Conference on Artificial Intelligence and Pattern Recognition*, pages 38–42, 2018. **3**
8. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. **2, 4, 8**
9. Jyh-Jing Hwang, Henrik Kretschmar, Joshua Manela, Sean Rafferty, Nicholas Armstrong-Crews, Tiffany Chen, and Dragomir Anguelov. Cramnet: Camera-radar fusion with ray-constrained cross-attention for robust 3d object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 388–405. Springer, 2022. **2, 3**

10. Youngseok Kim, Sanmin Kim, Jun Won Choi, and Dongsuk Kum. Craft: Camera-radar 3d object detection with spatio-contextual fusion transformer. *arXiv preprint arXiv:2209.06535*, 2022. [2](#), [3](#)
11. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [8](#)
12. Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2117–2125, 2017. [8](#)
13. Jianan Liu, Weiyi Xiong, Liping Bai, Yuxuan Xia, Tao Huang, Wanli Ouyang, and Bing Zhu. Deep instance segmentation with automotive radar detection points. *IEEE Transactions on Intelligent Vehicles*, 2022. [3](#)
14. Ramin Nabati and Hairong Qi. Centerfusion: Center-based radar and camera fusion for 3d object detection. In *Proceedings of the IEEE Workshop on Applications of Computer Vision (WACV)*, pages 1527–1536, 2021. [2](#), [3](#)
15. Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 652–660, 2017. [2](#)
16. Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Neural Information Processing Systems (NIPS)*, 30, 2017. [3](#)
17. Julien Rebut, Arthur Ouaknine, Waqas Malik, and Patrick Pérez. Raw high-definition radar for multi-task learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17021–17030, 2022. [1](#), [2](#), [4](#), [8](#), [9](#), [10](#), [11](#)
18. Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 115:211–252, 2015. [4](#)
19. Nicolas Scheiner, Florian Kraus, Nils Appenrodt, Jürgen Dickmann, and Bernhard Sick. Object detection for automotive radar point clouds—a comparison. *AI Perspectives*, 3(1):1–23, 2021. [3](#)
20. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [4](#)
21. Sourabh Vora, Alex H Lang, Bassam Helou, and Oscar Beijbom. Pointpainting: Sequential fusion for 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4604–4612, 2020. [1](#)
22. Chen Wang, Danfei Xu, Yuke Zhu, Roberto Martín-Martín, Cewu Lu, Li Fei-Fei, and Silvio Savarese. Densefusion: 6d object pose estimation by iterative dense fusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3343–3352, 2019. [5](#)
23. Bin Yang, Wenjie Luo, and Raquel Urtasun. Pixor: Real-time 3d object detection from point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7652–7660, 2018. [8](#), [9](#)
24. Ziguo Zhong, Stanley Liu, Manu Mathew, and Aish Dubey. Camera radar fusion for increased reliability in adas applications. *Electronic Imaging*, 2018(17):258–1, 2018. [3](#)